

# A Flexible, Scalable and Efficient Algorithmic Framework for *Primal* Graphical Lasso

Rahul Mazumder

Department of Statistics, Stanford University, Stanford, CA  
email: rahulm@stanford.edu

Deepak K. Agarwal

Yahoo! Research, 4401 Great America Parkway, Santa Clara.  
email: dagarwal@yahoo-inc.com

Submitted for publication on 10-11-2011

## Abstract

We propose a scalable, efficient and statistically motivated computational framework for Graphical Lasso (Friedman et al., 2007b) — a covariance regularization framework that has received significant attention in the statistics community over the past few years. Existing algorithms have trouble in scaling to dimensions larger than a thousand. Our proposal significantly enhances the state-of-the-art for such moderate sized problems and gracefully scales to larger problems where other algorithms become practically infeasible. This requires a few key new ideas. We operate on the primal problem and use a subtle variation of block-coordinate-methods which drastically reduces the computational complexity by orders of magnitude. We provide rigorous theoretical guarantees on the convergence and complexity of our algorithm and demonstrate the effectiveness of our proposal via experiments. We believe that our framework extends the applicability of Graphical Lasso to large-scale modern applications like bioinformatics, collaborative filtering and social networks, among others.

## keywords

Graphical Lasso,  $\ell_1$  regularization / LASSO, sparse inverse covariance selection, large scale convex optimization, convergence analysis, covariance estimation, positive definite matrices

## 1 Introduction

**Problem Description** Let  $\mathbf{S}_{p \times p}$  denote a  $p$ -dimensional sample covariance matrix obtained through i.i.d samples from a multivariate Gaussian distribution with (unknown) covariance  $\mathbf{\Sigma}$  and precision matrix  $\mathbf{\Sigma}^{-1}$ . The negative log-likelihood is given by:

$$f(\mathbf{\Theta}) := -\log \det \mathbf{\Theta} + \langle \mathbf{S}, \mathbf{\Theta} \rangle \text{ on } \mathbf{\Theta} \succ \mathbf{0}, \quad (1)$$

where  $\langle \mathbf{S}, \boldsymbol{\Theta} \rangle := \text{tr}(\mathbf{S}\boldsymbol{\Theta})$  and  $\boldsymbol{\Theta}$  corresponds to the precision matrix. The MLE (when it exists) is  $\hat{\boldsymbol{\Theta}} = \mathbf{S}^{-1}$ , but this estimator has high variance unless the sample size  $n$  is large relative to the dimension  $p$ . This makes the MLE a not-so-useful estimator of the covariance/precision matrix. In such high dimensional problems regularization (smoothing) is imperative to obtain reliable estimates. In fact, for the Gaussian distribution the precision matrix (Cox & Wermuth, 1996; Lauritzen, 1996) captures conditional dependencies among variables where absence of an edge (zero entry in the precision matrix) implies conditional independence. Hence, taking recourse to smoothing methods that induce sparsity is attractive. In addition to producing shrinkage estimators, a sparse precision graph leads to interpretable models and also provides model compression. In the context of learning large-scale graphs it is often undesirable from the point of view of computational/storage considerations to learn a model with all possible  $p^2$  edges present. Surprisingly enough, for large scale problems i.e. with  $p \approx 10^4 - 10^6$ , sparse precision graphs are computationally tractable, whereas their dense counterparts are not (Mazumder & Hastie, 2011, see for details.).

The  $\ell_1$  regularization (Friedman et al., 2007b; Banerjee et al., 2008; Yuan & Lin, 2007; Meinshausen & Bühlmann, 2006) is often used in this context since it performs smoothing, induces sparsity, adds interpretation and forms an effective model selection procedure. This is popularly known as *sparse inverse covariance selection* or the Graphical Lasso and is obtained as a solution to the following regularized criterion:

$$\underset{\boldsymbol{\Theta} \succ \mathbf{0}}{\text{minimize}} \quad g(\boldsymbol{\Theta}) := f(\boldsymbol{\Theta}) + \lambda \sum_{ij} |\theta_{ij}| \quad (2)$$

where  $\lambda > 0$  is the amount of regularization imposed on the entries of the precision matrix  $\boldsymbol{\Theta} = ((\theta_{ij}))$ . Equation (2) is a convex optimization problem (Semi-Definite Program aka SDP) in the variable  $\boldsymbol{\Theta}$ . The class of models described through equation (2) has already gained widespread interest in many statistical applications like biostatistics, functional magnetic resonance imaging, network analysis, collaborative filtering (Friedman et al., 2007b; Huang et al., 2010; Agarwal et al., 2011), and many more. Considerable progress has also been made in studying the statistical properties of the estimator and its variants (Ravikumar et al., 2011; Lam & Fan, 2009). We also note that the optimization in (2) is often used in a more non-parametric fashion (Agarwal et al., 2011; Mazumder & Hastie, 2011) for any positive semidefinite input matrix  $\mathbf{S}$ , not necessarily a sample covariance matrix from a MVN sample.

**Context and Background** Interior point methods for solving (2) scale poorly with increasing dimensions and quickly become infeasible for problem sizes around a hundred. For scalability purposes, first order methods relying on gradient information instead of Hessian (i.e. second order methods) become almost imperative (Nesterov, 2003). Over the past few years there has been substantial interest in developing such specialized scalable solvers for (2) (Friedman et al., 2007b; Banerjee et al., 2008; Lu, 2009, 2010; Scheinberg et al., 2010; Yuan, 2009; Boyd et al., 2011). However, existing solvers have difficulty in scaling to problems with  $p > 1000$  — precluding the wide-spread use of these methods in modern day

applications like collaborative filtering, graph mining, web-applications, large microarray data and other high dimensional problems.

There have been other interesting formulations to sparse precision matrix estimation (Rothman et al., 2010; Cai et al., 2011; Meinshausen & Bühlmann, 2006; Friedman et al., 2010, for example,). The formulation of Rothman et al. (2010) is non-convex. Cai et al. (2011) consider a linear programming approach where the precision matrix estimate *need not* be positive definite. Pseudo-likelihood based approaches (Friedman et al., 2010) do not ensure positive definiteness of the matrices. In this paper we focus on (2).

**Motivation** Several large scale covariance selection problems require algorithms that produce reasonably accurate approximations to the optimal solution of (2) within a certain time limit or equivalently, a limit on the computational budget (Bottou & Bousquet., 2008).

In fact, for large scale problems, under computational constraints an approximate solution to (2) is often the only feasible option. But for several applications, other than speed and scalability, it is necessary for the approximate solution to retain crucial properties of the optimal solution like sparsity and positive definiteness. One such application of large scale covariance selection was recently explored in the work of Agarwal et al. (2011). The authors used a sparse inverse covariance regularization to estimate the covariance matrix of high dimensional random effects in a multi-level hierarchical model. The paper explored prediction problems in recommender systems where the goal was to predict responses on unobserved user-item cells in a large matrix using responses on observed user-item pairs. Each user  $u$  is assigned an  $M$  dimensional random vector  $\phi_u$  that represents the user’s latent affinity to  $M$  items.  $\phi_u$ ’s are assumed to be drawn from a multivariate Gaussian prior with unknown covariance. The covariance is estimated via a E-M framework using an  $\ell_1$  regularization on the elements of precision matrix.

The use of a sparse inverse covariance regularization in the paper (Agarwal et al., 2011) led to a model with better predictive accuracy compared to other state-of-the-art methods. Since the estimation is based on an E-M strategy, it requires strictly positive definite estimates of the covariance and precision matrix. Indeed, an optimization of the form (2) needs to be conducted in the M-step — hence it is enough to terminate the process early without complete optimization. Early stopping along with sparsity in the precision matrix leads to a drastic reduction in computation time. The key property of covariance estimation required for the method to work is the ability to return both the precision matrix and its exact inverse i.e. the covariance matrix. These properties, apparently are not possessed by existing algorithms for (2). This may have precluded the use of sparse inverse covariance for estimating covariances in high dimensional random-effects model. We note that the strategy used in Agarwal et al. (2011) using the model fitting method described in this paper is general and can be used to model covariance in other high-dimensional multivariate random effects model that arise in applications like spatial statistics (Bernardinelli & Montomoli, 1992), social networks (Hoff, 2009; Hoff PD, 2002), and many more.

In the scenarios described above, we want a ‘flexible’ fitting algorithm for (2) such that:

- It can deliver a solution of arbitrary accuracy to (2) — the accuracy depending upon demands of the user/ application.

- Even if a low accuracy solution is desired, the algorithm upon exiting should return a sparse and positive definite  $\Theta$  and its inverse  $\Theta^{-1}$  — fundamental ingredients for relevant statistical model fitting procedures.
- The computational complexity per iteration of the algorithm is cheap enough to solve large scale problems.
- It readily adapts to warm-starts for computing a path of solutions on a grid of  $\lambda$  values.

We believe estimation procedures for inverse covariance described in this paper will make it routine to apply large scale covariance selection to high-dimensional multivariate data.

**Our Approach** We provide a brief outline of our approach and the salient features that make it different from other existing algorithms. Many of the sophisticated state-of-the-art algorithms (Banerjee et al., 2008; Lu, 2009, 2010; Scheinberg et al., 2010; Boyd et al., 2011) designed to solve (2) perform expensive operations like matrix inversions / eigen-decompositions on the *entire* matrix at every iteration requiring  $O(p^3)$ , which is clearly prohibitive for large problems. We take a different route by pursuing row/column block-coordinate based methods that cyclically update the estimates of one row/column at a time fixing the others at their latest values. Although Friedman et al. (2007b); Banerjee et al. (2008) also pursue block-coordinate methods, our approach differs in a few very important ways.

First, while we operate on the primal (where the primal variable is the precision matrix  $\Theta$ ), Friedman et al. (2007b); Banerjee et al. (2008) operate on the dual of (2). The primal and dual problems have some subtle and important differences that need consideration for large scale statistical applications. Algorithms operating on the dual (Friedman et al., 2007b; Lu, 2010; Banerjee et al., 2008, for example) do not return a sparse and positive definite precision matrix unless optimization is done to a very high degree of accuracy — this may be prohibitive for large scale problems. A more detailed discussion of this issue is provided in Section 6.

Second, we track both the precision and the covariance matrix over iterations, and our row/column block-coordinate wise procedure *does not perform a complete minimization* over the partial problems. This is a crucial observation since it reduces the row/column update cost from  $O(p^3)$  to  $O(p^2)$ . Although the idea looks simple at first blush, such incomplete minimization over partial problems is not necessarily guaranteed to ensure a proper optimization algorithm with convergence certificates. In Section 3 we show that such a relaxation still guarantees convergence of our algorithm. To the best of our knowledge, such a convergence analysis is novel both in the statistics and optimization literature.

**Contributions** We provide a summary of our main contributions before a detailed description of our approach. We propose a novel model fitting algorithm for a popular covariance selection method (2) that outperforms previous state of the art fitting algorithms for large problems with dimension  $p \approx 1 - 3$  thousands. The Algorithm design requires new and novel ideas. Our proposal is particularly suited to compute a path of solutions



to (2) by using warm-starts on a grid of  $\lambda$  values. The performance gains are quite impressive when compared to other existing algorithms for the same task as illustrated in Section 7. In addition, our algorithm is amenable to early stopping, provides a sparse and positive-definite solution, and scales to very large scale problems that are impractical to fit using existing methods. We provide a novel proof of asymptotic (algorithmic) convergence analysis, analyze complexity of the method and show the superiority of our methods through large scale simulation and data analysis. Finally, we outline how our approach can accommodate other row/column separable convex regularizers.

## 2 Algorithmic Framework

We now provide a detailed development of our fitting algorithm in this section, including convergence proof and computational complexity analysis. We begin with notations.

**Notations** We denote the set of all  $k \times k$  positive definite (respectively, positive semi-definite) matrices by  $S_k^{++}$  (respectively  $S_k^+$ ). We will write  $A_{k \times k} \succ 0$  if  $A \in S_k^{++}$ , similarly  $A \succeq 0$  implies  $A \in S_k^+$ . For a matrix  $A_{p \times p}$  we will denote its entries by  $a_{ij}, i = 1, \dots, p; j = 1, \dots, p$ .

For a vector  $\mathbf{u}$ , the notation  $\|\mathbf{u}\|_2$  denotes the usual  $\ell_2$  norm,  $\|\mathbf{u}\|_1$  denotes the  $\ell_1$  norm. For a matrix  $\mathbf{U}$ , we will use  $\|\mathbf{U}\|_2$  to denote its spectral norm i.e. the largest singular value of  $\mathbf{U}$ .

**Description of the Algorithm** The block coordinate method operates by fixing a row/-column index  $i \in \{1, 2, \dots, p\}$ , which without loss of generality, is assumed to be  $p$ . Partition the precision matrix  $\Theta$  and the sample covariance matrix  $\mathbf{S}$  as follows:

$$\Theta = \begin{pmatrix} \Theta_{11} & \boldsymbol{\theta}_{1p} \\ \boldsymbol{\theta}_{p1} & \theta_{pp} \end{pmatrix}, \quad \mathbf{S} = \begin{pmatrix} \mathbf{S}_{11} & \mathbf{s}_{1p} \\ \mathbf{s}_{p1} & s_{pp} \end{pmatrix}. \quad (3)$$

Using standard formulae for determinants of block-partitioned matrices we have:

$$\log \det(\Theta) = \log \det(\Theta_{11}) + \log(\theta_{pp} - \boldsymbol{\theta}_{1p}'(\Theta_{11})^{-1}\boldsymbol{\theta}_{1p}). \quad (4)$$

Using the above, the part of  $g(\Theta)$  in equation (2) that depends upon the  $p^{\text{th}}$  row/column of  $\Theta$  is given by:

$$g_p(\theta_{pp}, \boldsymbol{\theta}_{1p}) = -\log(\theta_{pp} - \boldsymbol{\theta}_{1p}'(\Theta_{11})^{-1}\boldsymbol{\theta}_{1p}) + 2\mathbf{s}_{1p}'\boldsymbol{\theta}_{1p} + (s_{pp} + \lambda)\theta_{pp} + 2\lambda\|\boldsymbol{\theta}_{1p}\|_1. \quad (5)$$

Note that the positive-definiteness of  $\Theta$  assures  $\theta_{pp} \geq 0$ . In (5), the optimization variables are  $\theta_{pp}$  and  $\boldsymbol{\theta}_{1p}$ . Conventional forms of block coordinate descent (Tseng, 2001; Friedman et al., 2007a) when applied to this problem will require completely minimizing the function (5) over the variables  $\boldsymbol{\theta}_{1p}$  and  $\theta_{pp}$ . Clearly for large problems an accurate optimization of this problem is quite computationally intensive, especially since this needs to be done several times across all rows/columns. We choose to deviate from this approach and propose to perform an *inexact minimization* in the afore-mentioned stage. The fact

that such a deviation still ensures a proper optimization procedure will be discussed later, for now we continue with the description of the algorithm.

Minimizing the criterion (5) with respect to  $\theta_{pp}$  with other coordinates fixed gives:

$$\hat{\theta}_{pp} := \arg \min_{\theta_{pp}} g(\theta_{pp}, \boldsymbol{\theta}_{1p}) = 1/(s_{pp} + \lambda) + \boldsymbol{\theta}'_{1p}(\boldsymbol{\Theta}_{11})^{-1}\boldsymbol{\theta}_{1p}. \quad (6)$$

The partially minimized objective (5), w.r.t.  $\theta_{pp}$  is given by:

$$\min_{\theta_{pp}} g_p(\theta_{pp}, \boldsymbol{\theta}_{1p}) = \log \det(s_{pp} + \lambda) + 2\mathbf{s}'_{1p}\boldsymbol{\theta}_{1p} + 1 + (s_{pp} + \lambda)\boldsymbol{\theta}'_{1p}(\boldsymbol{\Theta}_{11})^{-1}\boldsymbol{\theta}_{1p} + 2\lambda\|\boldsymbol{\theta}_{1p}\|_1.$$

Ignoring the constants independent of  $\boldsymbol{\theta}_{1p}$  above, we obtain an  $\ell_1$  regularized quadratic<sup>1</sup> function, which we denote by:

$$g_p(\boldsymbol{\theta}_{1p}) = \boldsymbol{\theta}'_{1p}\{(s_{pp} + \lambda)\boldsymbol{\Theta}_{11}^{-1}\}\boldsymbol{\theta}_{1p} + 2\mathbf{s}'_{1p}\boldsymbol{\theta}_{1p} + 2\lambda\|\boldsymbol{\theta}_{1p}\|_1. \quad (7)$$

We propose to use *one* sweep of cyclical coordinate-descent on this function  $g(\boldsymbol{\theta}_{1p})$ , w.r.t. the variable  $\boldsymbol{\theta}_{1p}$ .

We now summarize the update rule described above. Fix an arbitrary  $\tilde{\boldsymbol{\Theta}} \succ 0$ ,

$$\tilde{\boldsymbol{\Theta}} = \begin{pmatrix} \tilde{\boldsymbol{\Theta}}_{11} & \tilde{\boldsymbol{\theta}}_{1p} \\ \tilde{\boldsymbol{\theta}}_{p1} & \tilde{\theta}_{pp} \end{pmatrix} \quad (8)$$

and consider an increment in  $\boldsymbol{\Theta}$  around  $\tilde{\boldsymbol{\Theta}}$  in the direction of the  $p^{\text{th}}$  row/column. This updates  $\tilde{\boldsymbol{\Theta}}$  to  $\hat{\boldsymbol{\Theta}}$ :

$$\hat{\boldsymbol{\Theta}} \leftarrow \tilde{\boldsymbol{\Theta}} + (\boldsymbol{\omega}\mathbf{e}'_p + \mathbf{e}_p\boldsymbol{\omega}') \quad (9)$$

where  $\mathbf{e}_p$  is a vector in  $\Re^p$ , with all entries equal to 0 but the  $p^{\text{th}}$  entry equals to 1,  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_p)$  denotes the “increment” in the direction of the  $p^{\text{th}}$  row/column. Using

---

**Algorithm 1** Inner Block Inexact Coordinate-Descent

---

1. Initial value of  $\boldsymbol{\Theta}$  is  $\tilde{\boldsymbol{\Theta}}$ . Assign  $\hat{\boldsymbol{\Theta}} = \tilde{\boldsymbol{\Theta}}$ .
  2. Update the entries  $\omega_1, \dots, \omega_{p-1}$  and also  $\tilde{\boldsymbol{\theta}}_{1p}$ , as in (10) and (11).
  3. Update  $\hat{\omega}_p$  using the update-rule (13). Consequently change the  $(p, p)^{\text{th}}$  entry of  $\hat{\boldsymbol{\Theta}}$  to  $\hat{\omega}_p + \tilde{\theta}_{pp}$ .
- 

notation (9) and  $g_p(\cdot), \cdot \in \Re^{p-1}$  as in (7), the update rule in  $\boldsymbol{\omega}$  is given by:

$$\hat{\omega}_i = \arg \min_{\omega_i} g_p(\dots, (\hat{\boldsymbol{\theta}}_{1p})_{i-1}, (\tilde{\boldsymbol{\theta}}_{1p})_i + \omega_i, (\tilde{\boldsymbol{\theta}}_{1p})_{i+1}, \dots) \quad (10)$$

$$(\hat{\boldsymbol{\theta}}_{1p})_i \leftarrow (\tilde{\boldsymbol{\theta}}_{1p})_i + \hat{\omega}_i, \quad (\hat{\boldsymbol{\theta}}_{p1})_i \leftarrow (\tilde{\boldsymbol{\theta}}_{p1})_i + \hat{\omega}_i, \quad i = 1, \dots, p-1. \quad (11)$$

---

<sup>1</sup>note that the problem is convex only if  $\boldsymbol{\Theta}_{11}^{-1} \succeq 0$ , which is the case by virtue of the positive definiteness of the precision matrices, as shown in Section 2.1

Observe that the update (10) is simply a soft-thresholding operation:

$$\widehat{\omega}_i = \text{sgn}(a_i)(|a_i| - \lambda)_+ / b_i - (\widetilde{\boldsymbol{\theta}}_{1p})_i, \quad \text{where,} \quad (12)$$

$$a_i = (\mathbf{s}_{1p})_i + (s_{pp} + \lambda) \left( \sum_{j < i} (\boldsymbol{\Theta}_{11}^{-1})_{ij} (\widehat{\boldsymbol{\theta}}_{1p})_j + \sum_{j > i} (\boldsymbol{\Theta}_{11}^{-1})_{ij} (\widetilde{\boldsymbol{\theta}}_{1p})_j \right), \quad b_i = (s_{pp} + \lambda) (\boldsymbol{\Theta}_{11}^{-1})_{ii}$$

Finally, upon updating the off-diagonal entries in  $\widetilde{\boldsymbol{\Theta}}$ , the diagonal entry is updated using (6):

$$\widehat{\omega}_p \leftarrow 1/(s_{pp} + \lambda) + \widetilde{\boldsymbol{\theta}}'_{1p} (\widetilde{\boldsymbol{\Theta}}_{11})^{-1} \widetilde{\boldsymbol{\theta}}_{1p} - \widetilde{\theta}_{pp} \quad (13)$$

Overall, the above steps lead to the update formula:  $\widehat{\boldsymbol{\Theta}} \leftarrow \widetilde{\boldsymbol{\Theta}} + (\widehat{\boldsymbol{\omega}} \mathbf{e}'_p + \mathbf{e}_p \widehat{\boldsymbol{\omega}}')$ .

Note that the above operations require evaluations of the residual  $a_i$ . This requires computing at the onset the full gradient vector of the smooth part in (7) at  $\boldsymbol{\theta}_{1p}$ . When a coordinate of the vector  $\widetilde{\boldsymbol{\theta}}_{1p}$  gets updated, the entire gradient vector changes — this update can be achieved in  $O(p)$  flops. Note that in case  $\widehat{\omega}_i = 0$ , no update is required. Hence, if on *average* the number of non-zeros in  $\widehat{\boldsymbol{\theta}}_{p1}$ ,  $\widetilde{\boldsymbol{\theta}}_{p1}$  is  $k$ , then the update (10)–(11) requires an overall  $O(pk)$  flops which, for  $k \ll p$  leads to a significant reduction over the cost of a dense matrix/vector multiplication i.e.  $O(p^2)$ . Algorithm 1 summarizes the updating steps described above.

The above description was for updating the  $p^{\text{th}}$  row/column of the matrix  $\boldsymbol{\Theta}$ . This needs to be done for every row/column — one full sweep across the  $p$  rows/columns defines one iteration of our algorithm. We now describe the full version of our algorithm in Algorithm 2, we name it: **Primal INexact Minimization for Graphical Lasso (PINE-GL)**.

---

**Algorithm 2** Primal Inexact Minimization for Graphical Lasso (PINE-GL)

---

Inputs:  $\mathbf{S}, \lambda$ . Initialization:  $(\widetilde{\boldsymbol{\Theta}}, \widetilde{\boldsymbol{\Theta}}^{-1})$ .

- 1 For every row/column  $i \in \{1, 2, \dots, p, 1, 2, \dots\}$ , perform steps 2-3 till convergence.
  - 2 Permute the matrix such that the  $i^{\text{th}}$  row/column is the  $p^{\text{th}}$  i.e. of the form (3).  
 Obtain the matrix  $(\widetilde{\boldsymbol{\Theta}}_{11})^{-1}$  via rank-one updates (see Section 2.1.2).  
 Update the matrix using Steps 1 - 3 in Algorithm 1 :  $\widehat{\boldsymbol{\Theta}} \leftarrow \widetilde{\boldsymbol{\Theta}} + (\widehat{\boldsymbol{\omega}} \mathbf{e}'_p + \mathbf{e}_p \widehat{\boldsymbol{\omega}}')$   
 Obtain  $(\widehat{\boldsymbol{\Theta}})^{-1}$  via rank-one-updates (see Section 2.1.2).  
 Re-permute the matrix to get back the original rows/column indexing.
  - 3 Assign  $(\widetilde{\boldsymbol{\Theta}}, (\widetilde{\boldsymbol{\Theta}})^{-1}) \leftarrow (\widehat{\boldsymbol{\Theta}}, (\widehat{\boldsymbol{\Theta}})^{-1})$
  - 4 Upon convergence, the estimates at  $\lambda$ :  $(\widehat{\boldsymbol{\Theta}}_\lambda, (\widehat{\boldsymbol{\Theta}}_\lambda)^{-1}) := (\widehat{\boldsymbol{\Theta}}, (\widehat{\boldsymbol{\Theta}})^{-1})$
- 

**Convergence criterion** The convergence criterion is based upon the relative difference in objective values between two successive iterations (i.e. sweeps across all the  $p$  rows/columns), being less than a threshold. As described later, the objective value is computed on the ‘fly’, so expensive log-det computations need not be done separately.

**Initialization of precision and covariance matrices** PINE-GL requires as input, initialization for the tuple  $(\tilde{\Theta}, \tilde{\Theta}^{-1})$ . In case no prior choice for the input initialization is available, we use  $\tilde{\Theta}^{-1} \leftarrow \text{diag}(s_{11} + \rho, \dots, s_{pp} + \rho)$ . Note that the diagonals of  $\Theta^{-1}$  at the KKT optimality conditions for (2) is precisely the vector  $(s_{11} + \rho, \dots, s_{pp} + \rho)$ .

Often PINE-GL is used for computing a path of solutions to (2) via warm-starts — in such a case the tuple  $(\tilde{\Theta}, \tilde{\Theta}^{-1})$  is available as a by-product of the algorithm (see Section 2.2).

## 2.1 Important Properties of PINE-GL

We outline some of the important properties of our Algorithm — which is instrumental in making it flexible. For ease of exposition the technical details are relegated to the Supplementary Materials Section C.1.

### 2.1.1 Positive definiteness of precision & covariance matrices across the iterates

If the starting matrix  $\tilde{\Theta} \succ \mathbf{0}$ , then every row/column update in Step-2 of Algorithm 2 preserves positive definiteness of the updated matrix. For a rigorous proof see Section C.1.1 (supplementary materials).

### 2.1.2 Tracking precision and covariance matrices at every iteration

The function (7) that arises while updating the  $p^{\text{th}}$  row/column requires knowledge of  $(\tilde{\Theta}_{11})^{-1}$ . Of course, it is not desirable to compute the inverse from scratch for every row/column  $i$ , with a complexity of  $O(p^3)$ . However, if *both* the current precision and covariance matrices i.e.  $(\tilde{\Theta}, (\tilde{\Theta})^{-1})$  are stored at every iteration then it is quite simple to obtain  $(\tilde{\Theta}_{11})^{-1}$  via a matrix rank-one update as described in (32). This costs  $O(p^2)$  and moreover is amenable to parallelism. Similarly after every row/column update in  $\tilde{\Theta}$  its inverse can be obtained via a rank-one update as described in (33). Details of this implementation involve careful attention to details that are presented in the Section C.1.2 (supplementary materials).

Tracking both  $\Theta, \Theta^{-1}$  along the iteration provides flexibility to our algorithm in terms of:

- We avoid the additional cost of matrix inversion —  $O(p^3)$ .
- Termination at a given computational budget which is crucial for large scale problems and often desirable for exploratory analysis. Since the operating variable is  $\Theta$  — the precision matrix estimate is sparse <sup>2</sup>.
- They provide the perfect recipe for warm-starts, when one is interested in computing a path of solutions to (2) (see Section 2.2).

---

<sup>2</sup>This is different from the dual optimization problem, where the estimated positive definite precision matrices need not be sparse

- It gives a simple but efficient way to evaluate the log-determinant of the precision matrices along iterations, since computing the log-likelihood in large problems is a fairly expensive task.

## 2.2 Path Seeking Strategy

In many real life applications it is desirable to compute a path of solutions  $\{\hat{\Theta}_\lambda\}_\lambda$  over a grid of  $\lambda$ -values  $\lambda_K > \lambda_{K-1} > \dots > \lambda_1$ . One method is to compute the solutions across different tuning parameter values independently of each other, say on different machines. Otherwise, they can be computed serially wherein warm-starts/continuation strategies turn out to be very effective (Friedman et al., 2007a). In such a case, the estimate at  $\lambda_k$  i.e.  $(\hat{\Theta}_{\lambda_k}, (\hat{\Theta}_{\lambda_k})^{-1})$  is taken as an input<sup>3</sup> for the Algorithm 2 at  $\lambda = \lambda_{k-1}$ , for every  $k = K, \dots, 2$ . See Section 7 for experimental studies showing impressive improvements.

## 3 Convergence analysis

In this section we will analyze the convergence properties of Algorithm 2. We summarize below the novelty and importance of addressing convergence analysis in this paper:

Firstly, our proposal is not a conventional form of block coordinate descent as described in Tseng (2001); Friedman et al. (2007b), where the partial optimization problem (with the other variables fixed) is completely optimized. A complete-block coordinate minimization when applied to our problem requires a full minimization in Step 2 of Algorithm 2, over the  $i^{\text{th}}$  row/column. We differ by replacing this conventional *full* optimization strategy by a *partial* optimization — namely one pass of coordinate descent as described in Algorithm 1.

Secondly, our objective function of interest is non-smooth and due to the symmetry of the problem, the blocks i.e. the rows and columns have shared elements. Since  $\theta_{12} = \theta_{21}$  the value gets updated twice — once at row/column=1 and the other at row/column = 2. Conventional forms of block coordinate minimization theorems (Tseng, 2001) for non-smooth functions demand separability (in blocks) — so they do not apply directly. WEN et al. (2009) highlight this issue of overlapping entries and provide a proof of convergence. The work of WEN et al. (2009) considers smooth objectives — hence the results do not directly apply to our problem.

Note that by construction the sequence of precision matrices produced by Algorithm 2 results in a monotone decreasing sequence of objective values. Even if the objective values converge (which is true if they are bounded from below), it is not necessary that the point of convergence corresponds to the minimum of the problem (2) — the sequence of precision matrices need not converge either (see Tseng (2001) for discussions). We address these issues and show that the precision matrix estimates *converge* to the minimum under the mild assumption  $\lambda > 0$ . Convergence holds for  $\lambda = 0$  under the extra assumption that  $\mathbf{S} \succ \mathbf{0}$ .

The convergence analysis we present here is to the best of our knowledge novel.

---

<sup>3</sup>Note that the primal formulation is unconstrained — so warm-starts do not run into infeasibility problems.

We start with an important Lemma appearing in Lu (2010)[Proposition 3.1]:

**Lemma 1.** *For every  $\lambda > 0$ , problem (2) has a unique minimizer —  $\Theta_\lambda^*$ , which is (strictly) positive-definite and satisfies:*

$$\beta I_{p \times p} \geq \Theta_\lambda^* \geq \alpha I_{p \times p}$$

*for scalars  $\alpha, \beta$ , depending upon  $\mathbf{S}, \lambda, p$  with  $\infty > \beta \geq \alpha > 0$ .*

We are now ready to state the main theorem establishing the convergence of Algorithm 2

**Theorem 1** (Asymptotic Convergence of PINE-GL ). *Take  $\lambda > 0$ . Let  $\Theta_k$  be the estimate of the precision matrix obtained at iteration  $k$  i.e. on completion of Step 3 of Algorithm 2. Then the following hold true:*

(a) *The sequence of objective values is monotone decreasing :*

$$g(\Theta_{k+1}) \leq g(\Theta_k), \forall k \geq 1. \quad (14)$$

*The sequence  $\{g(\Theta_k)\}_k$  converges to the optimal solution of problem (2).*

(b) *The iterates  $\Theta_k \succ \mathbf{0}, \forall k$  and the sequence converges to  $\Theta_\infty$  — the unique solution to (2).*

*Proof.* The proof, which is rather detailed and technical is provided in the Appendix, Section A.1.  $\square$

## 4 Some variants of PINE-GL

This section discusses some variations to our proposal PINE-GL — leading to two important variants.

A variant of Algorithm 1 is having a counter for the number iterations for Step 2, say  $T_o$ . Our proposal of *inexact* minimization and for that matter the overall complexity analysis demands  $T_o = O(1)$  i.e.  $T_o \ll p$ . In our numerical experiments we found  $T_o \leq 2$  to be quite practical.

If  $T_o$  is taken to be arbitrarily large, we get the conventional form of cyclical coordinate descent used for  $\ell_1$  regularized quadratic programs (QP) (Friedman et al., 2007a). The magnitude of  $T_o$  depends upon the accuracy of the solution for the  $\ell_1$  regularized QP. In general, for a high accuracy solution, this can be arbitrarily large. If  $T_o = O(p)$  this leads to a  $O(p^3)$  complexity of Algorithm 1. See Section 5 for details. We call this variant **Primal Exact Minimization for Graphical Lasso** i.e. PEX-GL — this is the more conventional form of block coordinate descent applied on the problem (2).

We now proceed to discuss another simple but important variant of our algorithm PINE-GL, namely a ‘growing’ strategy — which we call **Primal GRowth for Graphical Lasso** i.e. PGR-GL.

## 4.1 Primal Growth for Graphical Lasso (PGR-GL )

Given an initial working dimension  $p_0$  (typically  $p_0 = 1$ ) and estimates of the precision and the covariance matrix  $(\tilde{\Theta}_{p_0 \times p_0}, (\tilde{\Theta}_{p_0 \times p_0})^{-1})$ , Algorithm 3 (Supplementary Materials, Section C.2) describes the task of obtaining the solution to (2) (with  $\mathbf{S}$  having dimension  $p \times p$ ). The main idea is to perform an initial forward pass by *incrementally* appending rows/columns and operating Step 2 of Algorithm 2 on the just added row/ column. Once the growing matrix is saturated to have  $p$  rows/columns — we make further passes through the  $p$  rows/columns, via Step-2 of Algorithm 2, till convergence. Since the ‘growing’ stage of the algorithm performs mainly cheap computations, it helps in providing pretty accurate warm-starts/ initializations  $\tilde{\Theta}, (\tilde{\Theta})^{-1}$  to PINE-GL , within a very short amount of time. See also results in Section 7. When the task is to solve (2) for a single value of  $\lambda$ , the method PGR-GL often turns out to be quite competitive with PINE-GL .

**Remark 1.** *The convergence of PGR-GL is straightforward. Firstly, it is not hard to see that the iterates maintain positive definiteness of the precision and its inverse and furthermore, since PINE-GL comes into action after one full-sweep of incrementally growing rows/columns, the convergence analysis for PINE-GL carries over.*

## 5 Computational Complexity

Here we describe the computational complexities of our proposed algorithms PINE-GL , PEX-GL and PGR-GL . We provide a summarized report here, the details are available in Appendix, Section B.

Cost of PINE-GL : Every row/column update requires  $O(p^2)$ , and for a full sweep across  $p$  rows/columns — this is  $O(p^3)$ . For  $\kappa$  full sweeps across  $p$  rows/columns this is  $O(\kappa p^3)$ , typically convergence occurs within  $\kappa \approx 2 - 10$ . See Section B.1.

Cost of PEX-GL : For every row/column the cost at the worst is  $O(p^3)$ . For a total of  $\kappa'$  ( $\approx 4-10$ ) sweeps across all rows/columns the cost is  $O(\kappa' p^4)$ . The cost may reduce to  $O(p^3)$  in case  $\lambda$  is quite large. See Section B.2.

Cost of PGR-GL : The cost here is  $O(p^3)$  as in PINE-GL — but the constants are generally better than that of PINE-GL . See Section B.3.

## 6 Related work

In this section we briefly describe some of the state-of-the art algorithms for criterion (2), their computational complexities and their fundamental differences with our proposal(s).

The block coordinate proposals of Banerjee et al. (2008); Friedman et al. (2007b) are related to our proposal — they solve the dual of the problem (2), which is given by:

$$\max_{\|\mathbf{V}\|_\infty \leq \lambda} -\log \det(\mathbf{S} + \mathbf{V}) - p. \quad (15)$$

By strong duality the optimal solution of problem (15) and (2) are the same, the primal-dual relationship being  $(\Theta)^{-1} = \mathbf{S} + \mathbf{V}$ . (15) operates on the covariance matrix whereas

the primal problem (2) operates on the precision matrix. As pointed out earlier, there is significant difference in pursuing the primal approach versus the dual. Often in real-life applications (as is the case in a principal motivating application for this paper (Agarwal et al., 2011)) one desires an approximate solution since it gives a fairly good statistical estimate for the main statistical estimation problem. An approximate solution in the dual space need not translate to one of similar accuracy in the primal space according to criterion (2). Further the dual approach does not produce sparse precision matrices — if  $\hat{\mathbf{V}}$  solves (15), then the precision matrix  $\hat{\Theta} = (\mathbf{S} + \hat{\mathbf{V}})^{-1}$  is *not* sparse unless (15) is solved till high tolerance ( $10^{-8}$ – $10^{-10}$ ). Ad-hoc thresholding strategies / post-processing strategies can be used to sparsify  $\hat{\Theta}$  — but positive definiteness is not guaranteed.

The block coordinate maximization of Banerjee et al. (2008) on (15), requires *solving* a box-constrained QP completely — with cost  $O(p^3)$ . The GLASSO (Graphical Lasso) Algorithm of Friedman et al. (2007b) minimizes the dual of the same box-constrained QP — an  $\ell_1$  regularized QP via cyclical coordinate descent. In the worst case this can be  $O(p^3)$ , in case the solutions are very sparse this is  $O(p^2)$ . Inexact minimization strategies for GLASSO do not guarantee convergence. GLASSO need not produce a sparse and positive definite precision matrix unless it converges to a high accuracy.

To summarize, both the block-coordinate proposals of Banerjee et al. (2008); Friedman et al. (2007b) have a worst case cost  $O(p^4)$  — the latter can improve to  $O(p^3)$  if  $\lambda$  is very large.

The gradient based algorithm of Banerjee et al. (2008) inspired by Nesterov (2005) has a per-iteration complexity  $O(p^3)$  and overall complexity  $O(\frac{p^{4.5}}{\epsilon})$  (for an  $\epsilon > 0$  accurate solution).

Another very efficient gradient-based algorithm is SMACS proposed in Lu (2010), which also solves the dual formulation. This has per iteration complexity  $O(p^3)$  (due to expensive matrix operations like eigen-decompositions, matrix inversions) and an overall complexity of  $O(\frac{p^4}{\sqrt{\epsilon}})$ .

The number of iterations taken by GLASSO (Friedman et al., 2007b; Banerjee et al., 2008) and PINE-GL (and its variants like PEX-GL, PGR-GL) i.e. full sweeps across all rows and columns are relatively small in most examples — of the order of 4-10. For a solution of similar accuracy, the number of gradient iterations for SMACS is often of the order of hundreds (or even more than a thousand) for problems of size 1000/1500.

It appears that most existing algorithms for solving the sparse covariance selection problem have a complexity of  $O(p^4)$  or possibly larger, depending upon the algorithm used and the desired accuracy of the solution — making computations for (2) almost impractical for values of  $p$  much larger than 1000/1500.

In contrast, every row/column update of PINE-GL is  $O(p^2)$  — overall for  $\kappa$  sweeps across all rows/columns this is  $O(\kappa p^3)$ , where  $\kappa$  denotes the total number of sweeps across all the rows/columns (See Section 5). This is clearly an order of magnitude improvement over existing algorithms and is further substantiated by our experiments.



## 7 Experimental Studies : synthetic examples

This section provides a comparison of our proposed fitting methods with some state-of-the-art algorithms for the optimization problem (2).

We use our main proposal PINE-GL , its close cousin PGR-GL , and the variant PEX-GL for comparisons.

Among the existing algorithms, Lu (2010) was observed to be better than the proposal of Lu (2009), so we used the former for our comparisons. Scheinberg et al. (2010); Yuan (2009); Boyd et al. (2011) describe algorithms based on the Alternating Direction Methods of Multipliers — among them the algorithm of Boyd et al. (2011) was publicly available at Stephen Boyd’s website. We experimented with this algorithm, but found it to be slower than GLASSO , so we did not include it for our comparisons.

We thus compared our proposals with two very efficient algorithms :

**GLASSO :** The algorithm of Friedman et al. (2007b). We used the MATLAB wrapper around their Fortran code — available at <http://www-stat.stanford.edu/~tibs/glasso/index.html>.

**SMACS :** denotes the algorithm of Lu (2010). We used the MATLAB implementation `smooth_covsel` available at [http://people.math.sfu.ca/~zhaosong/Codes/SMOOTH\\_COVSEL/](http://people.math.sfu.ca/~zhaosong/Codes/SMOOTH_COVSEL/).

Note that the default convergence criteria for the above two algorithms are different — GLASSO checks the successive changes in the diagonals of the precision matrix, whereas SMACS relies on duality gap. Moreover GLASSO and SMACS operate on the dual, whereas our proposals PINE-GL , PGR-GL , PEX-GL operate on the primal. Since, solving (2) is the main goal, to make comparisons fair, we compared the primal likelihoods of the estimates produced by the algorithms.

A relatively weak convergence criterion on the dual is typically quite far off from a sparse and positive definite precision matrix. The GLASSO algorithm tracks a precision matrix  $\Theta$  and covariance matrix  $\mathbf{W}$  along the iterations but  $(\Theta)^{-1} \neq \mathbf{W}$  and the discrepancy can be quite large before the algorithm converges to a high accuracy in the dual space. Furthermore, even if the estimated precision matrix (prior to convergence) is sparse it need not be positive definite. SMACS produces estimates of precision matrices  $\Theta$  along the iterations — though they are positive definite, they are dense. Arbitrary thresholding (to zero) of the smaller entries may destroy positive definiteness of the matrix.

Our proposal on the other hand at every iteration tracks the precision matrix (which is both sparse and positive definite) and its (exact) inverse.

In order to make the primal and dual problems comparable we consider the times taken by the algorithms to converge till a relatively high tolerance i.e.  $\text{TOL}=10^{-5}$ , where

$$\text{Convergence Test Criterion: } \frac{(g(\Theta_k) - \widehat{g(\Theta_*)})}{|\widehat{g(\Theta_*)}|} < \text{TOL}. \quad (16)$$

Here  $\Theta_k$  is the estimate of the precision matrix produced by the respective algorithm at the end of iteration  $k$ , and  $\widehat{g(\Theta_*)}$  is the estimate of the minimum of (2) obtained by taking the minimum over different algorithms after running them for a large number of iterations<sup>4</sup>.

---

<sup>4</sup>In our examples we ran PINE-GL , PEX-GL , PGR-GL (each) for 30 iterations. They were enough to give solutions till an accuracy of  $10^{-8}$ .

All of our computations are done in MATLAB 7.11.0 on a 3.3 GhZ Intel Xeon processor, with single-computational-thread computations enabled. Our codes are written in MATLAB and C <sup>5</sup>. GLASSO is written entirely in Fortran. SMACS is written in MATLAB — we don’t think this puts SMACS at a (timing) disadvantage, since the major computations are matrix operations which are very well optimized in MATLAB.

## 7.1 Algorithm Timings

The simulation examples we used were inspired by Lu (2010). The population precision matrix  $\Sigma_{p \times p}^{-1}$  having approximately 0.01 proportion of non-zeros is generated as follows. We generate a matrix  $A_{p \times p}$  with entries in  $\{-1, 0, 1\}$ , with proportion of non-zeros 0.01. The  $-1$  and  $1$  occur with equal probability.  $A$  is symmetrized via  $A \leftarrow 0.5 \cdot (A + A')$ . All the eigen-values of  $A$  are lifted up by adding a scalar multiple of a identity matrix :  $\Sigma^{-1} \leftarrow A + \tau I_{p \times p}$  such that the minimal eigen-value of  $\Sigma^{-1}$  is one. The (population) covariance matrix is taken to be  $\Sigma$ . We then generated  $x_i \sim MVN(0, \Sigma), i = 1, \dots, N$ . The sample correlation matrix was taken as  $\mathbf{S}$ .

We considered a battery of examples with varying  $N, p$ :

- (a)  $N \in \{500, 1000, 2000\}$  for  $p = 1000$  and (b)  $N \in \{2000, 3000, 4000\}$  for  $p = 1500$

Table 1 summarizes the timing results (in seconds) for the examples above for different algorithms. The timings are shown for different  $\lambda$  values — algorithms are cold-started at their default starting points. We see that PGR-GL, PINE-GL are always the winners and often by multiplicative factors. PINE-GL turns out to be the clear winner overall. PEX-GL turns out to be slower than PINE-GL — which supports our crucial idea of inexact minimization in the inner-blocks and also supports our complexity analysis. As expected, the timings for the block coordinate algorithms deteriorate for smaller values of  $\lambda$ . For dense problems (which are arguably harder problems for the primal formulation), PGR-GL consistently performs quite well. PEX-GL and GLASSO often perform similarly. SMACS tends to be quite slow for larger problems, when compared to the block coordinate counterparts. We found SMACS to be quite competitive for very small values of  $\lambda$  — but these (almost) unregularized solutions are not much statistically meaningful unless  $n > p$ . SMACS faced problems with convergence for  $n < p$  situations where  $\mathbf{S}$  was low-rank.

These results demonstrate the impressive comparative performances of PINE-GL and PGR-GL compared to current state-of-the-art methods — making it probably a method of choice in scenarios where it is possible to run the fitting algorithms till a high tolerance. As mentioned earlier, the primal formulation is particularly suited for this task of delivering solutions with lower tolerances. It operates on the primal (2) delivers a sparse and positive definite precision matrix and its exact inverse. Table 3 (Supplementary Material, Section C.3) shows average timings in seconds across a grid of ten  $\lambda$  values with varying degrees of accuracy. PINE-GL, PEX-GL and PGR-GL return lower accuracy solutions to (2) — in times which are much less than that taken to obtain higher accuracy solutions. The gains are rather substantial given the limited scope of the dual optimization problems in the ‘early stopping’ paradigm.

---

<sup>5</sup>The C code was generated via the embedded-matlab function `emlc`, an automated C code generator in Real Time Workshop in MATLAB.

p / N	% of nnz	Algorithm Times (sec)				
		PGR-GL	PEX-GL	PINE-GL	GLASSO	SMACS
$1 \times 10^3 / 2 \times 10^3$	92	<b>48.8</b>	140.4	119.4	143.8	308.5
	78	143.3	130.5	<b>94.7</b>	151.8	288.6
	46	149.7	167.4	<b>108.6</b>	220.1	217.7
$1 \times 10^3 / 1 \times 10^3$	85	79.1	143.1	<b>66.9</b>	130.4	398.7
	71	143.2	150.6	<b>69.0</b>	160.5	408.2
	49	<b>132.6</b>	225.3	187.2	295.7	464.2
$1 \times 10^3 / 0.5 \times 10^3$	91	82.4	93.8	<b>66.1</b>	94.4	—
	73	<b>81.5</b>	123.2	119.5	179.7	—
	48	354.8	382.4	<b>340.2</b>	544.5	—
$1.5 \times 10^3 / 4 \times 10^3$	86	223.2	258.7	<b>186.3</b>	571.1	2310.4
	72	221.8	353.4	<b>186.5</b>	577.8	1534.5
	47	<b>401.8</b>	656.1	488.4	851.8	1062.8
$1.5 \times 10^3 / 3 \times 10^3$	86	212.6	228.9	<b>177.4</b>	533.3	1736.3
	72	221.3	256.4	<b>186.0</b>	573.7	2017.2
	48	525.6	675.8	<b>494.2</b>	880.4	1521.4
$1.5 \times 10^3 / 2 \times 10^3$	85	283.3	364.8	<b>222.7</b>	566.5	1759.3
	72	222.6	258	<b>186.3</b>	574.2	2246.9
	40	757.8	1019.7	<b>706.6</b>	1186.6	1780.3

Table 1: Table showing the times in seconds till convergence to a tolerance of  $\text{TOL} = 10^{-5}$  (16), for different algorithms for different problem set-ups. For every combination of  $(p, N)$  three different  $r\lambda$  values are considered — as indicated by the % non-zeroes in the final solution for PEX-GL. All algorithms are warm-started at their default starting values. The ‘—’ corresponding to SMACS indicates that the algorithm did not converge for this example with  $N < p$ .

The next section compares different algorithms as ‘path-algorithms’. Path-based-algorithms and algorithms operating on a single value of  $\lambda$  are quite different performance-wise. An algorithm that tends to be very fast as a path algorithm need not be as good at a single value of  $\lambda$ . This is because, a good warm-start improves the convergence-rate of the algorithm. Similarly, an algorithm that is very good at a single value of  $\lambda$ , may not benefit much from warm-starts (a typical example being interior point methods). This is why we compare our proposals in both scenarios.

## 7.2 Tracing out a path of solutions

Continuing with Section 2.2, we see what happens to the rate of convergence of PINE-GL in presence of warm-starts. Note that SMACS and PGR-GL do not allow for warm-starts but GLASSO and PEX-GL do. The data is the same as used in the previous section. We took a grid of ten  $\lambda$  values as follows: the off-diagonal entries of the sample covariance matrix  $\mathbf{S}$  were sorted as per their absolute values and ten  $\lambda$  values were chosen from the entire range (along equi-spaced percentiles of the absolute values in  $\mathbf{S}$ ) — the largest  $\lambda$

value being  $\eta_{\max}q$  and the smallest was  $\eta_{\min} \cdot q$ , where  $q := \max_{i>j} |s_{ij}|$ . All algorithms were run till a tolerance of  $10^{-5}$ . Table 2 summarizes the results.

p / N	$\eta_{\max}/\eta_{\min}$ ( $\times 10^{-2}$ )	average % of nnz	Algorithm Times (sec)			speed-up (PINE-GL )
			PINE-GL	GLASSO	SMACS	
$1 \times 10^3 / 2 \times 10^3$	16/ 0.64	61.3	<b>77.27</b>	144.10	250.43	1.56
$1 \times 10^3 / 1 \times 10^3$	21/ 0.83	62.8	<b>116.38</b>	202.73	412.14	1.52
$1 \times 10^3 / 0.5 \times 10^3$	28 / 2	66.7	<b>105.52</b>	315.17	—	1.89
$1.5 \times 10^3 / 4 \times 10^3$	13.1 / 0.4	62.4	<b>260.54</b>	579.52	145.35	1.28
$1.5 \times 10^3 / 3 \times 10^3$	14.3 / 0.53	62.8	<b>280.19</b>	613.67	1631.6	1.23
$1.5 \times 10^3 / 2 \times 10^3$	16.0 / 0.61	63.1	<b>267.03</b>	697.65	1892.1	1.53

Table 2: Table showing the comparative timings (in seconds) of the three algorithms PINE-GL , GLASSO and SMACS across a grid of ten  $\lambda$  values. Times are *averaged* across the ten  $\lambda$  values. The averaged % of non-zeros in the final solution across the different  $\lambda$  values along with the limits of the  $\lambda$  values are also shown. The last column shows the speed-up factor for PINE-GL using warm-starts over the time spent to compute the solutions of the same accuracy without using warm-starts.

As the column on speed-up factor shows, the path algorithm of PINE-GL is much faster than obtaining the solutions at the same values of  $\lambda$  without warm-starts. PINE-GL continues to perform very well when compared to the path algorithm GLASSO .

## 8 Real Application: Learning Precision graphs for Movie-Movie Similarities

**MovieLens Data Set:** We study an application of the inverse covariance estimation method on a dataset obtained from a movie recommendation problem. We use the benchmark MovieLens-1M dataset available at <http://www.grouplens.org/node/12>, which consists of 1M explicit movie ratings by 6,040 users to 3,706 movies. The explicit ratings are on a 5-point ordinal scale, higher values indicative of stronger user preference for the movie. The statistical problem that has received considerable attention in the literature is that of predicting explicit ratings for missing user-movie pairs. Past studies (Salakhutdinov & Mnih, 2008) have shown that using movie-movie similarities based on “who-rated-what” information is strongly correlated with how users explicitly rate movies. Thus using such information as user covariates helps in improving predictions for explicit ratings. Other than using it as covariates, one can also derive a movie graph where edge weights represents movie similarities that are based on global “who-rated-what” matrix. Imposing sparsity on such a graph is attractive since it is intuitive that a movie is generally related to only a few other movies. This can be achieved through PINE-GL . Such a graph provides a neighborhood structure that can also help directly in better predicting explicit ratings. For instance, in predicting explicit rating  $r_{ij}$  by user  $i$  on movie  $j$ , one can use a weighted average of ratings by the user in the neighborhood of  $j$  derived from the movie-movie graph. Such neighborhood information can also be used as a graph Laplacian

to obtain better regularization of user factors in matrix factorization model as shown in Lu et al. (2009).

Other than providing useful information to predict explicit ratings, we note that using who-rated-what information also provides information to study the relationships among movies based on user ratings. We focus on such an exploratory analysis here but note that the output can also be used for prediction problems following strategies discussed above. A complete exploration of such strategies for prediction purposes in movie recommender applications is involved and beyond the scope of this paper.

We define the sample movie-movie similarity matrix as follows: for a movie  $j$ ,  $\mathbf{x}_j$  is the binary indicator vector denoting users who rated movie  $j$ . The similarity between movie  $j$  and  $k$  is defined as  $s_{jk} = \frac{\mathbf{x}_j' \mathbf{x}_k}{\sqrt{\sum x_{j,l} \sum x_{k,l}}}$ . The movie-movie similarity matrix  $\mathbf{S}$  thus obtained is symmetric and positive semi-definite.

## 8.1 Timing comparisons

As noted earlier, for this application we use criterion (2) in a non-parametric fashion, where our intention is to learn the connectivity matrix corresponding to the sparse inverse covariance matrix. We will first show timing comparisons of our method PINE-GL (the path-seeking version) along with GLASSO and SMACS. We see that PINE-GL is the only method that delivers solutions within a reasonable amount of time — the estimated precision matrices are used to learn the connectivity structure among the items.

We ran GLASSO for nine  $\lambda$  values — which were equi-spaced quantiles between the 8-th and 75-th percentile of the entries  $\{|s_{ij}|\}_{i>j}$  — this range also covers estimated precision matrices that are quite dense.

The path versions of PINE-GL, GLASSO and SMACS were used to obtain solutions on the chosen grid of  $\lambda$  values. The timings are summarized below:

- PINE-GL produced a path of solutions across the nine  $\lambda$  values using warm-starts in a *total* of 6.722 hours.
- The path version of GLASSO on the other hand, could not complete the same task of computing solutions to (2) on the same set of nine  $\lambda$  values, within two full days.
- We also tried to use SMACS for this problem, but it took more than 14 hours to compute the solution corresponding to a single value of  $\lambda$ .

The timing advantages should not come as a surprise given the computational complexity of PINE-GL is order(s) of magnitude better than its competitors. The performance gap becomes more prominent with increasing dimensions — traces of evidence were observed in Section 7.

## 8.2 Description of the Results

Figure 2 (Section C.4 of Supplementary Materials) displays the nature of the edge-structures and how they evolve across varying strengths of the shrinkage parameter  $\lambda$  for the *whole* precision-graphs produced by PINE-GL.

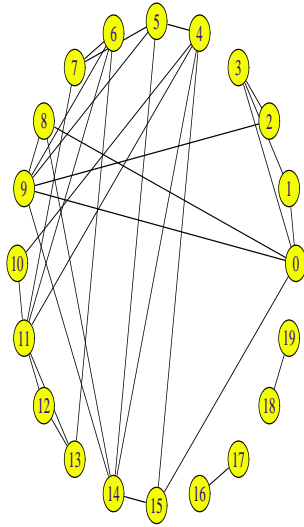
For a fixed precision matrix, a natural sub-set of ‘interesting’ edges among the all-possible  $p(p-1)/2$  edges are the ones corresponding to the top  $K$  absolute values of partial correlation coefficients. The nodes corresponding to these  $K$  edges and the edges of the precision graph restricted to them form a sub-graph of the  $p \times p$  precision graph. We summed the absolute values of the off-diagonal entries of the precision matrices across the different  $\lambda$  values. The (averaged)  $p(p-1)/2$  values were ordered and the top ten entries were chosen. These represent the partial correlations having the maximal strength (on average) across the different  $\lambda$  values taken. There were 20 vertices corresponding to these top ten partial correlation coefficients. Figure 1 shows the sub-graphs of the movie-movie precision graphs restricted to the selected 20 movies, across different  $\lambda$  values. The movie to ID mappings are in Table 4, in Section C.5 at Supplementary Materials Section. The edges in these subgraphs provide some interesting insights. For instance, consider the sub-graph corresponding to the largest  $\lambda$  (highest sparsity). Part of strong connectivity among movies 0,1,2,3 is expected since 0,1 and 2,3 are sequels. It is interesting to see there is a connection between 0 and 3, both of which are Sci-fi movies related to aliens. Other connections also reveal interesting patterns, these can be investigated using the IMDB movie database.

Another very related application of the set-up described above is the one appearing in Agarwal et al. (2011), where one models the raw who-rated-whom binary data using a multivariate random effects model.

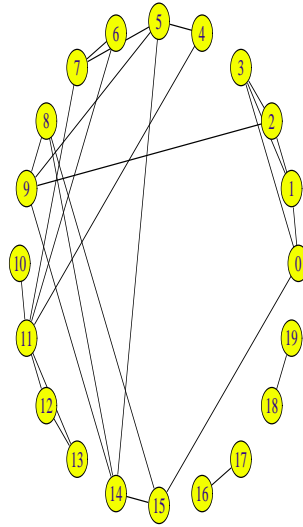
## 9 Conclusion and Remarks

We propose a flexible, scalable and efficient algorithmic framework for large scale  $\ell_1$  penalized inverse covariance selection problems that is used in several statistical applications. The framework gives rise to our main proposal PINE-GL, its close cousin PGR-GL and PEX-GL — all of them operate on the primal version of the problem (2). The key ingredient to scalability and efficiency requires a novel idea — that of inexact-minimization over an *exact*-one in the row/column blocks. The non-smoothness in the objective, positive definiteness of the precision matrices and the overlapping entries of the rows/columns necessitates a separate convergence analysis. We address this issue. This observation immediately brings down the per-iteration complexity of the algorithm by an order of magnitude, from  $O(p^3)$  to  $O(p^2)$ . On problems of size  $p = 1 - 3$  K, our proposal performs extremely well when compared to state of the art methods designed for problem (2). Our proposal tracks a sparse, positive definite precision matrix and its exact inverse i.e the covariance matrix at every iteration and is suited to return a solution with low/moderate accuracy depending upon the application task at hand. In particular, this makes it particularly suitable for large scale covariance selection problems where a very high accuracy solution is not practically feasible. PINE-GL is particularly suitable for computing a path of solutions on a grid of  $\lambda$  values using warm-starts — and it performs better when compared to existing path algorithms.

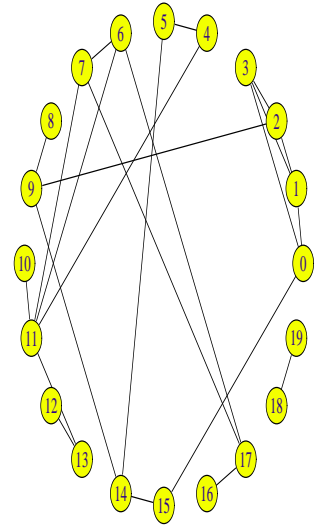
Our algorithm is supported by complexity analysis which shows that it is favorable over existing algorithms by an order of magnitude. Our proposals are well supported by numerical experiments on real and synthetic data.



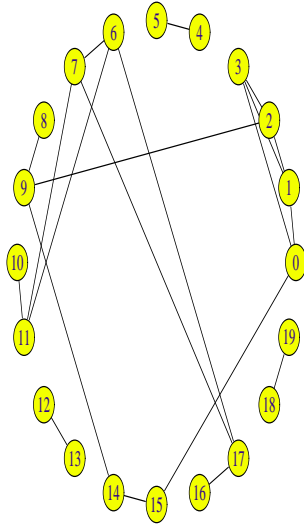
% non-zeros 90.4



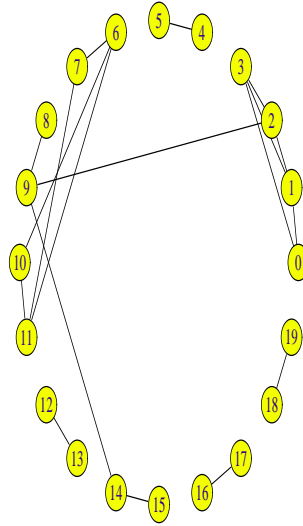
% non-zeros 95.0



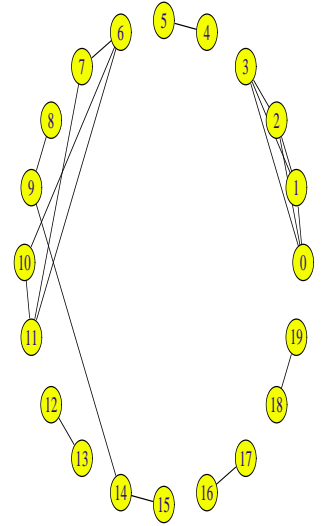
% non-zeros 97.8



% non-zeros 98.3



% non-zeros 98.8



% non-zeros 98.9

Figure 1: Pictures of subgraphs of the precision matrices induced by the 20 movies corresponding to the largest absolute partial correlations (averaged across different  $\lambda$  values).

**Exact Thresholding of Covariance Matrices** Fairly recently Mazumder & Hastie (2011) proposed an exact thresholding strategy which becomes useful if the non-zeros of the solution  $\Theta_\lambda^*$  to (2) breaks down into connected components. The *same* components can be recovered by looking at the non-zeros of the matrix  $\eta_\lambda(\mathbf{S})$ , where  $\eta(\cdot) = \text{sgn}(\cdot)(|\cdot| - \lambda)_+$  is the component-wise soft-thresholding operator at  $\lambda$ . As shown in Mazumder & Hastie (2011), this strategy can be used as a wrapper around any algorithm for solving (2) for sufficiently large  $\lambda$  so that it admits a decomposition into connected components. Since the aforementioned strategy heavily relies on having a scalable algorithm for (2), determined by the size of the maximal component — our proposal opens the possibility of solving problems (2) for an even wider range of  $\lambda$ -values.

**Extensions to other convex regularizers** Though we were concerned with (2) in this paper, our framework can accommodate other variants of block-separable regularizers, in place of the  $\ell_1$  norm on the entries of the matrix. This includes:

1. The weighted  $\ell_1$  norm i.e.  $\sum_{ij} p_{ij} |\theta_{ij}|$ , where  $p_{ij} \geq 0, \forall i, j$  are given scalars. See Friedman et al. (2007b); Fan et al. (2009) for use of this penalty.
2. The group lasso /node-sparse (Friedman et al., 2010) norm on the blocks of the precision matrix:  $\sum_{i=1}^p \sqrt{\sum_{j \neq i} \theta_{ij}^2}$ .
3. The elastic net regularization i.e.  $\alpha \sum_{ij} |\theta_{ij}| + (1 - \alpha) \sum_{ij} \theta_{ij}^2$  (Zou & Hastie, 2005)

These all are achieved by modifying Algorithm 1, with an inexact minimization strategy for the above penalties.

## 10 Acknowledgements

We would like to thank Trevor Hastie (Stanford University) and Liang Zhang (Yahoo! Labs) for helpful discussions. Many thanks to Liang for setting up the Movie-Lens data-set for our experiments. Rahul Mazumder would like to thank Yahoo ! Research for hospitality during the summer of 2010, during which this work started.



# A APPENDIX : Proofs and Technical Details

This Section accumulates the technical details and proof details that were outlined in the text.

## A.1 Theorem 1 : Set up and Proof

Firstly we will like to point out some important points about the convergence result which also sheds important light on the properties of the solution (2). If  $\lambda > 0$ , the sequence of objective values and the estimates are bounded below (see Lemma 1). Then by standard results in real-analysis, the sequence of objective values converge to  $g_\infty$  (say). It is not clear however (see Tseng (2001), and references therein for counter-examples) that the point of convergence i.e.  $g_\infty$  corresponds to the minimum of the problem (2). Fortunately however, we show in this section that  $g_\infty$  actually is the optimum of the minimization problem (2).

Observe that the convex optimization problem (2), for  $\lambda = 0$  and  $\mathbf{S}$  rank-deficient will have its infimum at  $-\infty$ . However, it turns out that for  $\lambda > 0$ , this condition is avoided and the optimal value of the problem is finite.

As is the case for many convex optimization problems (Boyd & Vandenberghe, 2004, see for example), its is not necessary that problem (2) will have a unique minimizer. It turns out, however, that as soon as  $\lambda > 0$ , problem (2) has a unique minimizer. The assertions made above are consequences of Lemma 1.

We need to set up a formal framework and present a few lemmas leading to the proof.

**The  $\ell_1$  Regularized Proximal Map** A variant of Step 2, in Algorithm 1 is one where we use a proximal step(Nesterov, 2007), instead of one sweep of cyclical block-coordinate descent. Recall that the function  $g_p(\boldsymbol{\theta}_{1p})$  (7) is in the composite form (Nesterov, 2007) :

$$g_p(\boldsymbol{\theta}_{1p}) = f_p(\tilde{\boldsymbol{\theta}}_{1p}) + 2\lambda\|\boldsymbol{\theta}_{1p}\|_1 \quad (17)$$

where  $f_p(\cdot)$  denotes the smooth part given by:

$$f_p(\tilde{\boldsymbol{\theta}}_{1p}) = \boldsymbol{\theta}'_{1p} \{ (s_{pp} + \lambda) \boldsymbol{\Theta}_{11}^{-1} \} \boldsymbol{\theta}_{1p} + 2\mathbf{s}'_{1p} \boldsymbol{\theta}_{1p}$$

It is easy to see that the gradient  $\nabla f_p(\cdot)$  of the function  $f_p(\cdot)$  is Lipschitz continuous i.e. :

$$\|\nabla f_p(\mathbf{u}) - \nabla f_p(\mathbf{v})\|_2 \leq L_p \|\mathbf{u} - \mathbf{v}\|_2 \quad (18)$$

and an estimate of  $L_p = 2(s_{pp} + \lambda)\|\boldsymbol{\Theta}_{11}^{-1}\|_2$ . The proximal step or the generalized gradient step (in place of the coordinate-wise updates (10) ) is given by the following:

$$\hat{\boldsymbol{\omega}} = \arg \min_{\boldsymbol{\omega} \in \mathbb{R}^{p-1}} \left\{ \frac{L_p}{2} \|\boldsymbol{\omega} - (\tilde{\boldsymbol{\theta}}_{1p} - \frac{1}{L_p} \nabla f_p(\tilde{\boldsymbol{\theta}}_{1p}))\|_2^2 + 2\lambda\|\boldsymbol{\omega}\|_1 \right\} - \tilde{\boldsymbol{\theta}}_{1p} \quad (19)$$

$$= \eta(\tilde{\boldsymbol{\theta}}_{1p} - \frac{1}{L_p} \nabla f_p(\tilde{\boldsymbol{\theta}}_{1p}); \frac{2\lambda}{L_p}) - \tilde{\boldsymbol{\theta}}_{1p} \quad (20)$$

where  $\nabla f_p(\tilde{\boldsymbol{\theta}}_{1p}) = 2(s_{pp} + \lambda)(\tilde{\boldsymbol{\Theta}}_{11})^{-1}\tilde{\boldsymbol{\theta}}_{1p} + 2\mathbf{s}_{1p}$  and  $\eta(\cdot; \gamma) = \text{sgn}(\cdot)(|\cdot| - \gamma)_+$  is the soft-thresholding operator applied component-wise to a vector  $\cdot \in \mathbb{R}^{p-1}$ .

In what follows, we will study a minor variation in Step 2 of Algorithm 1. Instead of using one sweep of cyclical coordinate descent, we will use one proximal step as described in (20). The convergence result with the cyclical coordinate-descent version is no different but simply adds to the technicality of the analysis.

**Properties of the soft-thresholding operator** Before going into the proof we need a lemma pertaining to an important property of the soft-thresholding operator i.e. the  $\ell_1$  Regularized Proximal Map.

For a function  $h : \mathfrak{R}^q \mapsto \mathfrak{R}$  with Lipschitz continuous gradient:

$$\|\nabla h(\mathbf{x}) - \nabla h(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\| \quad (21)$$

the following majorization property holds (Beck & Teboulle, 2009, See for example, Lemma 2.1)

$$\frac{L}{2}\|\mathbf{w} - \mathbf{x}\|_2^2 + \langle \nabla h(\mathbf{x}), \mathbf{w} - \mathbf{x} \rangle + h(\mathbf{x}) \geq h(\mathbf{w}) \quad (22)$$

The minimizer wrt  $\mathbf{w}$  for the  $\ell_1$  regularized problem:

$$\text{Maj}(\mathbf{w}; \mathbf{x}) := \frac{L}{2}\|\mathbf{w} - \mathbf{x}\|_2^2 + \langle \nabla h(\mathbf{x}), \mathbf{w} - \mathbf{x} \rangle + h(\mathbf{x}) + \lambda\|\mathbf{w}\|_1 \quad (23)$$

is given by the proximal map or the soft-thresholding operator:

$$\begin{aligned} M(\mathbf{x}) &:= \arg \min_{\mathbf{w}} \left\{ \frac{L}{2}\|\mathbf{w} - \mathbf{x}\|_2^2 + \langle \nabla h(\mathbf{x}), \mathbf{w} - \mathbf{x} \rangle + \lambda\|\mathbf{w}\|_1 \right\} \\ &= \arg \min_{\mathbf{w}} \left\{ \frac{L}{2}\|\mathbf{w} - (\mathbf{x} - \frac{1}{L}\nabla h(\mathbf{x}))\|_2^2 + \lambda\|\mathbf{w}\|_1 \right\} \\ &= \eta \left( (\mathbf{x} - \frac{1}{L}\nabla h(\mathbf{x})); \frac{\lambda}{L} \right) \end{aligned} \quad (24)$$

The following Lemma states an important property of the map  $M(\mathbf{x})$ .

**Lemma 2.** *Consider the function  $H(\cdot)$  defined by:*

$$H(\mathbf{w}) = h(\mathbf{w}) + \lambda\|\mathbf{w}\|_1 \quad (25)$$

*with  $h(\cdot)$  having the property in (21). For any  $\mathbf{x} \in \mathfrak{R}^q$  and  $M(\cdot)$  as defined in (24) the following holds:*

$$\frac{2}{L} \cdot (H(\mathbf{x}) - H(M(\mathbf{x}))) \geq \|\mathbf{x} - M(\mathbf{x})\|_2^2 \quad (26)$$

*Proof.* It can be shown using elementary convex analysis and the properties of the map  $\text{Maj}(\cdot)$  (23) defined above, (Beck & Teboulle, 2009, See Lemma 2.3):

$$H(\mathbf{x}) - H(M(\mathbf{y})) \geq \frac{L}{2}\|M(\mathbf{y}) - \mathbf{y}\|_2^2 + L\langle \mathbf{y} - \mathbf{x}, M(\mathbf{y}) - \mathbf{y} \rangle \quad (27)$$

Substituting  $\mathbf{y} = \mathbf{x}$  above we get the desired result in (26).  $\square$

**Proof of Theorem 1, part (a)** The monotonicity follows by construction of the sequence of iterates  $\Theta_k$ .

The iterate  $\Theta_k$  is obtained by updating all the  $p$  rows/columns of the matrix  $\Theta$ , cyclically. We now introduce some notation. Let us denote the successive row/column updates by:

$$\begin{aligned} \text{Update row/column 1} &\rightarrow \Theta_{k,1} \\ \text{Update row/column 2} &\rightarrow \Theta_{k,2} \\ &\dots\dots\dots \\ \text{Update row/column } p &\rightarrow \Theta_{k,p} \end{aligned} \tag{28}$$

Further we use  $\Theta_{k,i}[-i, i] \in \mathbb{R}^{p-1}$  to denote the  $i^{\text{th}}$  column of the matrix  $\Theta_{k,i}$  (excluding the diagonal entry). We need to estimate the difference in  $\Theta_{k,i-1}$  and  $\Theta_{k,i}$  — note that they differ only in the  $i^{\text{th}}$  row/column.

To settle ideas and using the same set-up as in Section 2, we concentrate on row/column  $p$ . The difference  $\Theta_{k,p}[-p, p] - \Theta_{k,p-1}[-p, p]$  can be estimated by using Lemma 2. To see this recall the framework of Algorithm 2 as described in Section 2. To update the  $p^{\text{th}}$  row/column we need to consider a proximal-gradient step in the function  $g_p(\cdot)$  as described in (17). This function exactly fits into the framework of Lemma 2, for specific choices of  $L$ ,  $h(\cdot)$ ,  $H(\cdot)$  and  $\lambda$ . Let the Lipschitz constant at this iterate be denoted by  $L_{k,p}$ . Using the equality  $g(\Theta_{k,p-1}) - g(\Theta_{k,p}) = g_p(\Theta_{k,p-1}[-p, p]) - g_p(\Theta_{k,p}[-p, p])$  (which follows by construction) and Lemma 2 we have:

$$\frac{2}{L_{k,p}}(g(\Theta_{k,p-1}) - g(\Theta_{k,p})) \geq \|\Theta_{k,p}[-p, p] - \Theta_{k,p-1}[-p, p]\|_2^2 \tag{29}$$

Recall that we established in Section 2.1.1, that Algorithm 2 produces a sequence of estimates  $\Theta_{k,i}$  such that  $\Theta_{k,i} \succ 0$ . Further note that the minimum of (2) is finite (as  $\lambda > 0$ , Lemma 1). It follows that there exists  $\rho' > \rho > 0$  such that

$$\rho' I_{p \times p} \succ \Theta_{k,i} \succ \rho I_{p \times p}, \forall k \tag{30}$$

where  $I_{p \times p}$  is a  $p$  dimensional identity matrix. Since  $L_{k,i}$  is a scalar multiple (18) of the spectral norm  $\|(\Theta_{k,i}[-i, -i])^{-1}\|_2$ , it follows from (30) that  $\inf_{k,i} L_{k,i} > 0$  and  $\infty > \sup_{k,i} L_{k,i}$ .

Thus using the monotonicity of the sequence of objective values  $g(\Theta_{k,i})$  for  $i = 1, \dots, p$ ,  $k \geq 1$ , the fact that the minimum value of (2) is finite and the boundedness of  $\frac{1}{L_{k,i}}$ , we see that the left hand side of (29) converges to zero as  $k \rightarrow \infty$ . This implies that  $\Theta_{k,p}[-p, p] - \Theta_{k,p-1}[-p, p] \rightarrow 0$  as  $k \rightarrow \infty$  i.e. the successive difference of the off-diagonal entries converge to zero as  $k \rightarrow \infty$ . In particular, we have this to be true for every row/column  $i \in \{1, \dots, p\}$  i.e.

$$\Theta_{k,i}[-i, i] - \Theta_{k,i-1}[-(i-1), i-1] \rightarrow 0, \quad k \rightarrow \infty \text{ for every } i \in \{1, \dots, p\}$$

In the above, we use the convention  $\Theta_{k,0}[-1, 1] = \Theta_{k-1,p}[-p, p]$ .

Since  $\{\Theta_k\}_k$  is a bounded sequence it has a limit point — let  $\Theta_\infty$  be a limit point. Moving along a sub-sequence (if necessary),  $n_k \subset \{1, 2, \dots\}$  we have  $\Theta_{n_k} \rightarrow \Theta_\infty$ .

Using the stationary condition for the update described in (20), the  $p^{\text{th}}$  row/column (off-diagonal entries) satisfies:

$$(s_{pp} + \lambda)(\Theta_\infty)_{11}^{-1}(\Theta_{\infty,p}[-p, p]) + \mathbf{s}_{1p} + \lambda \text{sgn}(\Theta_{\infty,p}[-p, p]) = 0.$$

The above holds true for every row/column  $i \in \{1, \dots, p\}$ . Using the above stationary condition along with the update relation for the diagonal entries as in Step 3 in Algorithm 1, it is easy to see that the limit point  $\Theta_\infty$  satisfies the global stationary condition for problem (2) i.e.

$$-(\Theta_\infty)^{-1} + \mathbf{S} + \lambda \text{sgn}(\Theta_\infty) = \mathbf{0}$$

Thus we have established that  $g(\Theta_k)$  converges to the global minimum of the function  $g(\cdot)$ .

**Proof of Theorem 1, part (b)** For this part it suffices to show that there is a unique limit point for the sequence  $\Theta_k$ . Note that we showed in Part (a) that every limit point of the sequence  $\Theta_k$  is a minimizer for the problem (2). Now by Lemma 1, there is a unique minimizer of (2). This implies that  $\Theta_k$  has a unique limit point and hence the sequence converges to  $\Theta_\infty$ , the minimizer of  $g(\cdot)$ .

## B Complexity analysis details of PINE-GL ,PEX-GL and PGR-GL

### B.1 Complexity of PINE-GL

Step 2 of Algorithm 2 requires computing  $\Theta_{11}$ , this requires  $O((p-1)^2)$  (see Section 2.1). Algorithm 1 does one sweep of cyclical coordinate descent — this has worst case complexity  $O((p-1)^2)$  — in case the solution to the  $\ell_1$  regularized QP is sparse, the cost is much smaller. It should be noted here that any constant number of cycles (say  $T_o$ ) of cyclical coordinate descent will lead to a complexity of  $O(T_o \cdot (p-1)^2)$ . As long as  $T_o \ll p$  (say  $T_o = 1/2$ ) this leads to  $O(p^2)$ . This is followed by updating the covariance matrix with  $O((p-1)^2)$ , via rank-one updates. Hence Step 2, for each row / column has a complexity of  $O(p^2)$ , for  $p$  rows/columns this is  $O(p^3)$ . If  $\kappa$  denotes the total number of full sweeps across all the rows and columns this leads to  $O(\kappa p^3)$ . In practice based on our experiments we found  $\kappa = 2 - 10$  is sufficient for convergence till a fairly high tolerance. While computing a path of solutions with warm-starts  $\kappa$  is around 2-4 for different values of  $\lambda$ . The value of  $\kappa$  increases when  $\lambda$  is very small so that the resultant solution  $\Theta^*$  is dense — but since these  $\lambda$  values correspond to almost un-regularized likelihood solutions, in most applications they are not statistically interesting solutions.

### B.2 Complexity of PEX-GL

In case of using PEX-GL the analysis is quite similar to above but there are subtle differences. The complexity of matrix rank-updates remain the same  $O(p^2)$ , what changes is the number of coordinate sweeps required for Algorithm 1 to solve the inner  $\ell_1$  regularized block QP till high accuracy. This problem is fairly challenging in its own right and is

computationally hard when  $p$  is a few thousand. Precise convergence rates of coordinate descent to the best of our knowledge are not known. This depends largely upon the data type being used. Often the number of coordinate sweeps i.e.  $k$  can be  $O(p)$  — leading to a complexity of  $O(p^3)$ . If (generalized) gradient descent methods (Nesterov, 2007) are used instead of cyclical coordinate descent — then the number of iterations  $k$  is of the order of  $O(1/\epsilon)$ , where  $\epsilon > 0$  is the accuracy of the solution. For  $\epsilon \approx 1/p$ ,  $k \approx p$ . Thus PEX-GL has roughly a complexity of  $O(p^3)$  for every row/column update — leading to an overall cost of  $O(p^4)$  for one full sweep across all rows/columns. If there are  $\kappa'$  full sweeps across all rows/columns the total cost is  $O(\kappa'p^4)$ .

In case  $\lambda$  is large enough so that every  $\ell_1$  regularized QP can be solved quite fast i.e.  $O(p^2)$  — the total cost of PEX-GL reduces to  $O(p^3)$ .

### B.3 Complexity of PGR-GL

The main difference of PGR-GL lies in the manner in which it updates the rows/columns via appending rows/columns in Steps 1-3 of Algorithm 3. Steps 1-3 have a cost of  $\sum_{m=1}^p m^2$ , which is approximately  $p^3/3$ . When compared with one sweep of PINE-GL, the growing step is approximately one-third of the cost of PINE-GL. One sweep of the growing strategy leads to inferior performance when compared to one sweep of PINE-GL. However, after a smaller number of sweeps, PGR-GL can obtain better likelihoods than PINE-GL. In some examples, as seen in the experimental section, PGR-GL is faster than PINE-GL in obtaining a solution with low accuracy.

## C SUPPLEMENTARY MATERIALS

This portion gathers some of the technicalities avoided in the main text of the article and the Appendix A.

### C.1 Properties of the updates of Algorithm 2

We present here a detailed derivation of the properties of the updates of Algorithm 2.

#### C.1.1 Positive-definiteness

The updates described in Steps (1), (2), (3) in Algorithm 1 actually preserve positive definiteness of  $\hat{\Theta}$ , under the assumption that  $\tilde{\Theta} \succ \mathbf{0}$ . Using the decomposition for  $\tilde{\Theta}$  in (3), it follows from standard properties of positive definiteness of block partitioned matrices (Boyd & Vandenberghe, 2004) that:

$$\tilde{\Theta} \succ \mathbf{0} \quad \text{iff} \quad \tilde{\Theta}_{11} \succ \mathbf{0}, \quad \tilde{\theta}_{pp} - \tilde{\theta}'_{1p} \tilde{\Theta}_{11}^{-1} \tilde{\theta}_{1p} > 0 \quad (31)$$

Observe that  $\hat{\Theta}_{11} = \tilde{\Theta}_{11}$ , by construction. Further by the property of the update Step 3 (Algorithm 2) we see that

$$\hat{\theta}_{pp} - \hat{\theta}'_{1p} \tilde{\Theta}_{11}^{-1} \hat{\theta}_{1p} = \frac{1}{s_{pp} + \lambda} > 0.$$

This shows by (31) that  $\hat{\Theta} \succ \mathbf{0}$ .

A simple consequence of the above observation is that  $\log \det(\hat{\Theta})$  is finite if  $\log \det(\tilde{\Theta})$  is so.

#### C.1.2 Tracking $(\hat{\Theta}, (\hat{\Theta})^{-1})$

For updating the  $p^{\text{th}}$  row/column, PINE-GL requires knowledge of  $(\tilde{\Theta}_{11})^{-1}$ . Of course, it is not desirable to compute the inverse from scratch for every row/column  $i$ , with a complexity of  $O(p^3)$ . Assume that, before operating on the  $p^{\text{th}}$  row/column we already have with us the tuple  $(\tilde{\Theta}, (\tilde{\Theta})^{-1})$  — then it is fairly easy to compute  $(\tilde{\Theta}_{11})^{-1}$  via:

$$(\tilde{\Theta}_{11})^{-1} = \tilde{\mathbf{W}}_{11} - \tilde{\mathbf{w}}_{1p} \tilde{\mathbf{w}}_{p1} / \tilde{w}_{pp}, \quad (32)$$

where  $\tilde{\mathbf{W}} := (\tilde{\Theta})^{-1}$  and the blocks  $\tilde{\mathbf{W}}_{11}, \tilde{\mathbf{w}}_{1p}, \tilde{\mathbf{w}}_{p1}, \tilde{w}_{pp}$  of the matrix  $\tilde{\mathbf{W}}$ , have the same structure as in (3). This follows by standard-formulae of inverses of block-partitioned matrices — and the update requires  $O(p^2)$ .

Once the  $p^{\text{th}}$  row/column of the matrix  $\tilde{\Theta}$  is updated, we obtain  $\hat{\Theta}$ . The matrix  $\hat{\mathbf{W}} := (\hat{\Theta})^{-1}$  is obtained via:

$$\hat{\mathbf{W}}_{11} = (\tilde{\Theta}_{11})^{-1} - \frac{(\tilde{\Theta}_{11})^{-1} \hat{\theta}_{1p} \hat{\theta}_{p1} (\tilde{\Theta}_{11})^{-1}}{(\hat{\theta}_{pp} - \hat{\theta}_{p1} (\tilde{\Theta}_{11})^{-1} \hat{\theta}_{1p})}; \quad \hat{\mathbf{w}}_{1p} = -\frac{(\tilde{\Theta}_{11})^{-1} \hat{\theta}_{1p}}{(\hat{\theta}_{pp} - \hat{\theta}_{p1} (\tilde{\Theta}_{11})^{-1} \hat{\theta}_{1p})}, \quad (33)$$

---

**Algorithm 3** PINE-GL with Growing Strategy : PGR-GL

---

Inputs:  $\lambda$ ,  $\mathbf{S}_{p \times p}$  and  $p_0 \times p_0$  matrices  $(\tilde{\Theta}, (\tilde{\Theta})^{-1})$ , where  $p_0 < p$ .

Set initial working row/column  $m = p_0 + 1$ .

- 1 Consider a  $m \times m$  dimensional problem of the form (2) with covariance<sup>6</sup>  $\mathbf{S}_{1:m \times 1:m}$  and initializations  $(\tilde{\Theta}, (\tilde{\Theta})^{-1})$ , of dimension  $m \times m$  where

$$\tilde{\Theta} \leftarrow \text{blkdiag}(\tilde{\Theta}, \frac{1}{(s_{mm} + \lambda)}); \quad (\tilde{\Theta})^{-1} \leftarrow \text{blkdiag}((\tilde{\Theta})^{-1}, s_{mm} + \lambda) \quad (34)$$

- 2 Apply Step 2 of Algorithm 2 with inputs  $\mathbf{S}_{1:m \times 1:m}$ ,  $(\tilde{\Theta}, (\tilde{\Theta})^{-1})$  and input dimension  $m$ .

The above results in  $(\hat{\Theta}, (\hat{\Theta})^{-1})$  — both  $m$  dimensional matrices.

Assign  $(\tilde{\Theta}, (\tilde{\Theta})^{-1}) \leftarrow (\hat{\Theta}, (\hat{\Theta})^{-1})$

- 3 Increase  $m = m + 1$ , and go to Step-1 (if  $m \leq p$ ); else go to Step-4.
  - 4 Apply Algorithm 2 with  $\mathbf{S}_{1:p \times 1:p}$ , and initializations  $(\tilde{\Theta}_{p \times p}, (\tilde{\Theta}_{p \times p})^{-1})$ , till a desired tolerance. The output is the solution to problem (2).
- 

where as before the blocks of the matrix  $\widehat{\mathbf{W}}$ , follow the same notation as in (3). The cost is again  $O(p^2)$ . Note that the diagonal entry  $\hat{w}_{pp} = 1/(\hat{\theta}_{pp} - \hat{\theta}_{p1}(\tilde{\Theta}_{11})^{-1}\hat{\theta}_{1p})$ .

The above recursion shows how to track  $(\hat{\Theta}, \hat{\Theta}^{-1})$  (as well as  $(\tilde{\Theta}, \tilde{\Theta}^{-1})$ ) as one cycles across the rows/columns of the matrix  $\Theta$ .

## C.2 Algorithmic Description of Primal Growth for Graphical Lasso PGR-GL

This elaborates Section 4.1, in the text. Given an initial working dimension  $p_0$  (typically  $p_0 = 1$ ) and estimates of the precision and the covariance matrix  $(\tilde{\Theta}_{p_0 \times p_0}, (\tilde{\Theta}_{p_0 \times p_0})^{-1})$ , Algorithm 3 describes the task of obtaining the solution to (2), with  $\Theta, \mathbf{S}$  having dimensions  $p \times p$ .

## C.3 Performances of PINE-GL and its variants for low–high accuracy solutions

We now proceed to show how gracefully they deliver solutions of lower accuracy within a much shorter span of time making it feasible to scale to very high dimensional problems. As mentioned earlier, the primal formulation is particularly suited for this task of delivering solutions with lower convergence tolerance, since it delivers a sparse and positive definite precision matrix and its exact inverse. Table 3 shows average timings in seconds across a grid of ten  $\lambda$  values with varying degrees of accuracy. In case the application demands a

relatively low accuracy solution, then the algorithms deliver solutions within fractions of the time taken to deliver a solution with higher accuracy.

p / N	average % of nnz	Accuracy (TOL)	Algorithm Times (sec)		
			PGR-GL	PEX-GL	PINE-GL
$1 \times 10^3 / 2 \times 10^3$	61.3	$10^{-2}$	60.16	68.29	<b>49.0</b>
		$10^{-3}$	80.59	104.67	<b>67.31</b>
		$10^{-4}$	112.31	134.13	<b>91.91</b>
		$10^{-5}$	140.08	174.52	<b>120.75</b>
$1 \times 10^3 / 1 \times 10^3$	62.77	$10^{-2}$	68.46	92.32	<b>64.56</b>
		$10^{-3}$	<b>95.49</b>	126.25	<b>95.47</b>
		$10^{-4}$	128.94	167.35	<b>127.11</b>
		$10^{-5}$	<b>174.44</b>	199.36	177.06
$1 \times 10^3 / .5 \times 10^3$	66.67	$10^{-2}$	82.94	117.97	<b>65.94</b>
		$10^{-3}$	127.17	153.29	<b>96.69</b>
		$10^{-4}$	181.99	195.19	<b>140.37</b>
		$10^{-5}$	252.95	234.97	<b>200.43</b>
$1.5 \times 10^3 / 4 \times 10^3$	62.44	$10^{-2}$	149.89	195.19	<b>124.55</b>
		$10^{-3}$	<b>189.52</b>	317.64	204.77
		$10^{-4}$	<b>266.63</b>	396.82	275.75
		$10^{-5}$	344.26	460.29	<b>333.55</b>
$1.5 \times 10^3 / 3 \times 10^3$	62.78	$10^{-2}$	145.59	215.65	<b>141.87</b>
		$10^{-3}$	203.86	300.57	<b>201.62</b>
		$10^{-4}$	<b>261.12</b>	397.72	271.34
		$10^{-5}$	<b>344.19</b>	477.64	346.47
$1.5 \times 10^3 / 2 \times 10^3$	63.11	$10^{-2}$	<b>149.13</b>	251.51	169.37
		$10^{-3}$	250.02	354.75	<b>238.57</b>
		$10^{-4}$	319.36	445.38	<b>317.51</b>
		$10^{-5}$	414.91	523.38	<b>408.77</b>

Table 3: Table showing average timings in seconds across a grid of ten  $\lambda$  values with varying degrees of accuracy i.e. TOL. The column under average % of non-zeroes denotes the % of non-zeroes in the optimal solution, averaged across the ten  $\lambda$  values. No warm-start across  $\lambda$ 's is used.

Table 3 shows that PINE-GL , PEX-GL and PGR-GL return lower accuracy solutions to (2) — in times which are much less than that taken to obtain higher accuracy solutions. Note that even the lower accuracy solutions correspond to sparse and positive definite precision matrices, with guarantees on ‘TOL’. Even more interesting is the flexibility of methods like PINE-GL , PGR-GL to obtain sparse and positive definite solutions at low computational cost when compared to the times taken by the dual algorithms in Table 1. These favorable comparative timings go on to support our claim of making large scale covariance selection very practical. PEX-GL turns out to be the slowest among the three, PGR-GL and PINE-GL are quite strong competitors, with PINE-GL winning in majority of the situations.



## C.4 Graphical display of sparsity patterns in the Movie-lens Graphs

This is an elaborate version of Section 8.2 in the main text. Figure 2 represents the sparsity structures of the precision graphs obtained from the movie-movie similarities. The graph structures are displayed under the Sparse reverse Cuthill-McKee ordering (Gilbert et al., 1992)<sup>7</sup> of the precision matrices as delivered by our algorithm PINE-GL for three different values of  $\lambda$ . The presence of a ‘dot’ in the figure represents a non-zero edge weight in the corresponding movie-movie precision graph. The percentages of (off-diagonal) non-zeros are presented below each figure. The tapering nature of the graphs for larger values of  $\lambda$ , show that the movies towards the extreme ends of the axes tend to be connected to few other movies. These movies tend to be conditionally dependent on very few other movies. The higher density of the points towards the center (of the left two figures) show that those movies tend to be connected to a larger number of other movies.

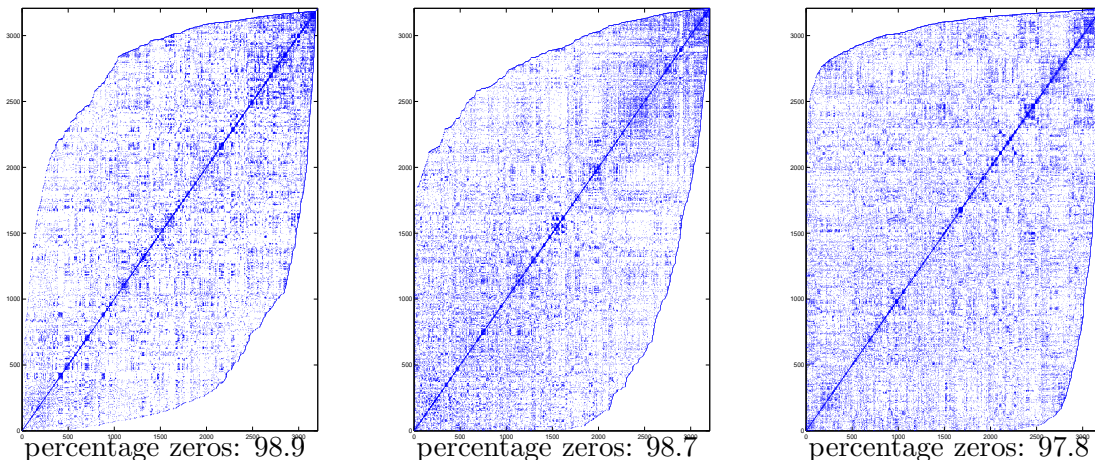


Figure 2: MATLAB `spy` plots under the reverse Sparse reverse Cuthill-McKee ordering of the vertices of the sparse precision matrices obtained via PINE-GL, for three different values of the tuning parameters. A dot represents presence of an edge. The percentage of off-diagonal zeros in the matrix are given below each plot.

## C.5 Movie-ID to Movie mapping Table

The movie ids— movie name mapping is given in the following table:

---

<sup>7</sup>For a sparse symmetric matrix  $A$  the reverse Cuthill-McKee ordering is a permutation  $\pi$  such that  $A(\pi, \pi)$  tends to have its nonzero elements closer to the diagonal. This is often used for visualizing the sparsity structure of large dimensional graphs

(0) PuppetMaster5: TheFinalChapter (1994)	(1) PuppetMaster4 (1993)
(2) Carnosaur3: PrimalSpecies (1996)	(3) Carnosaur2(1995)
(4) Fridaythe13thPartV: ANewBeginning (1985)	(5) Fridaythe13thPartVII:TheNewBlood(1988)
(6) Porky’sRevenge (1985)	(7) Porky’s II: TheNextDay (1983)
(8) SororityHouseMassacre (1986)	(9) SororityHouseMassacreII (1990)
(10) PoliceAcademy5: Assignment:MiamiBeach(1988)	(11) PoliceAcademy6:CityUnderSiege(1989)
(12)RockyIV(1985)	(13) RockyIII (1982)
(14) Hellbound:HellraiserII(1988)	(15) Hellraiser(1987)
(16) CloseShave,A(1995)	(17) WrongTrousers,The (1993)
(18) Godfather:PartII,The(1974)	(19) Godfather,The(1972)

Table 4: Table showing the names of the top 20 Movies, appearing in the top twenty strongest partial correlations. It is seen from Figure 1 that edges often occur between movies and their sequels.

## References

- AGARWAL, D., ZHANG, L. & MAZUMDER, R. (2011). Modeling item-item similarities for personalized recommendations on yahoo! front page. *Annals of Applied Statistics* **5**(3), 1839–1875.
- BANERJEE, O., GHAOUI, L. E. & D’ASPREMONT, A. (2008). Model selection through sparse maximum likelihood estimation for multivariate gaussian or binary data. *Journal of Machine Learning Research* **9**, 485–516.
- BECK, A. & TEOULLE, M. (2009). A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sciences* **2**, 183–202.
- BERNARDINELLI, L. & MONTOMOLI, C. (1992). Empirical bayes versus fully bayesian analysis of geographical variation in disease risk. *Statistics in Medicine* **11**, 9831007.
- BOTTOU, L. & BOUSQUET, O. (2008). The trade-offs of large scale learning. In *Advances in Neural Information Processing Systems 20*, J. Platt, D. Koller, Y. Singer & S. Roweis, eds. Cambridge, MA: MIT Press, pp. 161–168.
- BOYD, S., PARIKH, N., CHU, E., PELEATO, B. & ECKSTEIN, J. (2011). Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning* **3**(1), 1–122.
- BOYD, S. & VANDENBERGHE, L. (2004). *Convex Optimization*. Cambridge University Press.
- CAI, T., LIU, W. & LUO, X. (2011). A constrained l1 minimization approach to sparse precision matrix estimation. *Journal of the American Statistical Association* **106**, 594–607.
- COX, D. & WERMUTH, N. (1996). *Multivariate Dependencies*. Chapman and Hall, London.

- FAN, J., FENG, Y. & WU, Y. (2009). Network exploration via the adaptive lasso and scad penalties. *Annals of Applied Statistics* , 521–541.
- FRIEDMAN, J., HASTIE, T., HOEFLING, H. & TIBSHIRANI, R. (2007a). Pathwise coordinate optimization. *Annals of Applied Statistics* **2**, 302–332.
- FRIEDMAN, J., HASTIE, T. & TIBSHIRANI, R. (2007b). Sparse inverse covariance estimation with the graphical lasso. *Biostatistics* **9**, 432–441.
- FRIEDMAN, J., HASTIE, T. & TIBSHIRANI, R. (2010). Applications of the lasso and grouped lasso to the estimation of sparse graphical models.
- GILBERT, J. R., MOLER, C. & SCHREIBE, R. (1992). Sparse matrices in matlab: Design and implementation. *SIAM Journal on Matrix Analysis* .
- HOFF, P. D. (2009). Multiplicative latent factor models for description and prediction of social networks. *Comput. Math. Organ. Theory* **15**, 261–272.
- HOFF PD, RAFTERY AE, H. M. (2002). Latent space approaches to social network analysis. *Journal of the American Statistical Association* **97**, 1090–1098.
- HUANG, S., LI, J., SUN, L., YE, J., FLEISHER, A., WU, T., CHEN, K. & REIMAN., E. (2010). Learning brain connectivity of alzheimers disease by sparse inverse covariance estimation. *NeuroImage* **50**, 935–949.
- LAM, C. & FAN, J. (2009). Sparsistency and rates of convergence in large covariance matrix estimation. *Annals of Statistics* **37(6B)**, 4254–4278.
- LAURITZEN, S. (1996). *Graphical Models*. Oxford University Press.
- LU, Z. (2009). Smooth optimization approach for sparse covariance selection. *SIAM J. on Optimization* **19**, 1807–1827.
- LU, Z. (2010). Adaptive first-order methods for general sparse inverse covariance selection. *SIAM J. Matrix Anal. Appl.* **31**, 2000–2016.
- LU, Z., AGARWAL, D. & DHILLON, I. S. (2009). A spatio-temporal approach to collaborative filtering. In *RecSys*.
- MAZUMDER, R. & HASTIE, T. (2011). Exact covariance thresholding into connected components for large-scale graphical lasso. *arXiv* : 1108.3829v2, (*submitted*) .
- MEINSHAUSEN, N. & BÜHLMANN, P. (2006). High-dimensional graphs and variable selection with the lasso. *Annals of Statistics* **34**, 1436–1462.
- NESTEROV, Y. (2003). Introductory lectures on convex optimization: Basic course. *Kluwer, Boston* .
- NESTEROV, Y. (2005). Smooth minimization of non-smooth functions. *Math. Program., Serie A* **103**, 127–152.

- NESTEROV, Y. (2007). Gradient methods for minimizing composite objective function. Tech. rep., Center for Operations Research and Econometrics (CORE), Catholic University of Louvain. Tech. Rep, 76.
- RAVIKUMAR, P., RASKUTTI, G., WAINWRIGHT, M. J. & YU, B. (2011). High-dimensional covariance estimation by minimizing l1-penalized log-determinant. *Electronic Journal of Statistics(to appear)* .
- ROTHMAN, A. J., LEVINA, E. & ZHU, J. (2010). A new approach to cholesky-based covariance regularization in high dimensions. *Biometrika* **97**, 539–550.
- SALAKHUTDINOV, R. & MNIH, A. (2008). Probabilistic matrix factorization. In *The Twenty-Second Annual Conference on Neural Information Processing Systems*. MIT Press.
- SCHEINBERG, K., MA, S. & GOLDFARB, D. (2010). Sparse inverse covariance selection via alternating linearization methods. *NIPS* , 1–9.
- TSENG, P. (2001). Convergence of a block coordinate descent method for nondifferentiable minimization. *Journal of Optimization Theory and Applications* **109**, 475–494.
- WEN, Z., GOLDFARB, D. & SCHEINBERG, K. (2009). Row by row methods for semidefinite programming. *Industrial Engineering* , 1–21.
- YUAN, M. & LIN, Y. (2007). Model selection and estimation in the gaussian graphical model. *Biometrika* **94**, 19–35.
- YUAN, X. (2009). Alternating direction methods for sparse covariance selection. at:<http://www.optimization-online.org/DB-FILE/2009/09/2390.pdf> , 1–12.
- ZOU, H. & HASTIE, T. (2005). Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society Series B.* **67**, 301–320.